

1-1-1990

On the Optimal Reward Function of the Continuous Time Multiarmed Bandit Problem

José Luis Menaldi

Wayne State University, menaldi@wayne.edu

Maurice Robin

Institut National de Recherche en Informatique et en Automatique, Rocquencourt

Recommended Citation

J.-L. Menaldi and M. Robin, *On the optimal reward function of the continuous time multiarmed bandit problem*, SIAM J. Control Optim., 28 (1990), pp. 97-112. doi: [10.1137/0328005](https://doi.org/10.1137/0328005)

Available at: <http://digitalcommons.wayne.edu/mathfrp/35>

This Article is brought to you for free and open access by the Mathematics at DigitalCommons@WayneState. It has been accepted for inclusion in Mathematics Faculty Research Publications by an authorized administrator of DigitalCommons@WayneState.

ON THE OPTIMAL REWARD FUNCTION OF THE CONTINUOUS TIME MULTIARMED BANDIT PROBLEM*

JOSÉ LUIS MENALDI† AND MAURICE ROBIN‡

Abstract. The optimal reward function associated with the so-called “multiarmed bandit problem” for general Markov-Feller processes is considered. It is shown that this optimal reward function has a simple expression (product form) in terms of individual stopping problems, without any smoothness properties of the optimal reward function neither for the global problem nor for the individual stopping problems. Some results relative to a related problem with switching cost are obtained.

Key words. variational inequality, switching problem, bandit problem, dynamic programming, index policy

AMS(MOS) subject classifications. 35B37, 49A60, 49B60, 60J25, 93E20

1. Introduction. This paper deals with the properties of the optimal reward function associated with the so-called “multiarmed bandit problem.” Let us recall, formally, the statement of the problem: assume that there are N independent machines. $x_i(t)$, $t \in \mathbb{R}_+$ is the state (for instance the production) of machine i . At each time t , one operates only one machine, the others being frozen. When machine i is operating, $x_i(t)$ evolves as a continuous time Markov process with a given semigroup $\Phi^i(t)$. If $i(t)$ denotes the number of the machine in operation at time t , we want to maximize a global payoff

$$(1.1) \quad J = E \int_0^\infty e^{-\alpha t} f(i(t), x_{i(t)}(t)) dt$$

where f is a given instantaneous reward.

The multiarmed bandit problem has been studied by Gittins [4] and Whittle [8] in the discrete time case, and more recently by Varaiya, Walrand, and Buyukkoc [7] in a more general setting. Karatzas [5] studied the continuous time case when $x_i(t)$ is a one-dimensional diffusion process. The most general study is done in Mandelbaum [13], [14] who formulated the problem as the control of a multiparameter process. This approach allows, in particular, a strong formulation of the optimal process when $x_i(t)$ is a diffusion process.

In Whittle [9] it is shown that the optimal reward function has a simple expression in terms of an individual stopping problem each involving only one machine. Such an expression is shown to hold true for the diffusion bandit problem in Karatzas [5] thanks to the smoothness of the reward function which allows explicit computations.

In this paper, the main objective is to obtain such an expression when the $x_i(t)$ are general Feller processes, without smoothness properties of the optimal reward function neither for the global problem nor for the individual stopping problem.

Let us describe briefly what expression we are looking for.

* Received by the editors October 19, 1987; accepted for publication (in revised form) March 9, 1989.

† Wayne State University, Department of Mathematics, Detroit, Michigan 48202. This research was partially supported by National Science Foundation grant DMS-8601998 and Air Force Office of Scientific Research contract F49620-86-C-0111.

‡ Institut National de Recherche en Informatique et en Automatique, Rocquencourt, Domaine de Voluceau, B.P. 105, 78153 Le Chesnay, France.

Following Whittle [9], we will use the variant of the problem where one can decide, at any time, to stop the control problem, with a reward M if this “retirement option” is chosen.

Assume that $x_i(t)$ is for each i , a Markov process with values in some space E_i , with semigroup $\Phi^i(t)$.

If \underline{x} denotes the initial state of the whole set of machines, and if $u(\underline{x}, M)$ is the corresponding optimal reward function, then by applying, formally, the dynamic programming arguments, $u(\underline{x}, M)$ is shown to be the minimum solution of the following inequalities:

$$(1.2) \quad \begin{aligned} u(\underline{x}, M) &\geq e^{-\alpha t} \Phi^i(t) u(\underline{x}, M) + \int_0^t e^{-\alpha s} \Phi^i(s) f_i(x_i) ds \\ u(\underline{x}, M) &\geq M. \end{aligned}$$

The individual stopping problems have optimal cost functions $(\phi_i(x_i, M), i = 1, N)$, where ϕ_i is the minimum solution of

$$(1.3) \quad \begin{aligned} \phi_i(x_i, M) &\geq e^{-\alpha t} \Phi^i(t) \phi_i(x_i, M) + \int_0^t e^{-\alpha s} \Phi^i(s) f_i(x_i) ds \\ \phi_i(x_i, M) &\geq M \end{aligned}$$

when $\alpha k \leq f_i(x_i) \leq \alpha K, \forall i, \forall x_i$.

The objective is to show that

$$(1.4) \quad u(\underline{x}, M) = K - \int_M^K \prod_{i=1}^N \frac{\partial \phi_i}{\partial m} dm.$$

It would be nice to obtain such a formula by analytic methods, as it can be shown that (1.2) and (1.3) have a minimal solution (cf. [1], [2], [3]). However, without smoothness on ϕ_i , we do not know how to show the result by analytic methods.

Here we will use an intermediary control problem (§ 2.1) which is suitable for our objective, although it does not contain a general statement of the multiarmed bandit itself when there is no switching cost.

Using this particular interpretation of the minimal solution of (1.2), we will show (1.4) using an extension to the continuous time case of the Tsitsiklis’ lemma [6]. In § 3, we investigate the problem with switching cost, showing a similar lemma; it does not seem possible, however, to obtain an expression of the optimal reward in terms of some individual problems.

2. Problem without switching cost. We start with a control problem which will provide a stochastic interpretation of (1.2).

2.1. An intermediary control problem. Let $E_i, i = 1 \cdots N$ be a family of compact metric spaces endowed with their Borel σ -algebra.

Define $E = E_1 \times \cdots \times E_N$. Throughout the paper,

\underline{x} will denote an element of E , i.e.,

$$\underline{x} = (x_1, \cdots, x_N), \quad x_i \in E_i.$$

We are given a family of Markov semigroups $\Phi^i(t) i = 1, \cdots, N, \Phi^i(t)$ being defined and *continuous* on $C(E_i)$, the Banach space of continuous functions on E_i .¹

¹ So, Φ^i is a Feller semigroup on $C(E_i)$, cf. Dynkin [10].

If $\Omega_i = D(\mathbb{R}_+, E_i)$, the space of right continuous, left limited functions on \mathbb{R}_+ with values in E_i , we denote by $Q_{x_i}^i$ the probability measure on Ω_i corresponding to Φ^1 , and we define

$$\Omega = \Omega_1 \times \cdots \times \Omega_N$$

and $\{F_t\}$ the associated canonical σ -algebra.

In order to define the controlled process, we first consider the probability measure corresponding to constant trajectories for the components $j \neq i$ (i being the number of the process which is active, the others being frozen), and which gives the markovian evolution corresponding to $\Phi^1(t)$ for the component i : in other words we define

$$(2.1) \quad P_{i,x} = \delta_{x_1} \times \cdots \times \delta_{x_{i-1}} \times Q_{x_i}^i \times \delta_{x_{i+1}} \cdots \times \delta_{x_N}.$$

Notice that, if

$$\underline{x}_t(\omega) = \omega(t) \quad \text{for } \omega \in \Omega,$$

then

$$E_{i,x} g(\underline{x}_t) = E_{x_i}^i g(x_1, \cdots, x_{i-1}, x_i(t), x_{i+1}, \cdots, x_N)$$

where $E_{i,x}$ (respectively, $E_{x_i}^i$) denotes the expectation with respect to $P_{i,x}$, (respectively, $Q_{x_i}^i$).

Assume now that

$$(2.2) \quad f_i(x_i) \text{ is a positive function } f_i \in C(E_i), \quad \forall i \alpha > 0 \text{ a discount factor}$$

$$(2.3) \quad V \text{ will be the set of admissible controls and } v \in V \Leftrightarrow v = (\theta_n, \xi_n)_{n \geq 0}, \theta_0 = 0, \text{ where } (\theta_n) \text{ is an increasing sequence of } F_t \text{ stopping times, } \xi_n \text{ a } F_{\theta_n}\text{-measurable random variable with values in } \{1, \cdots, N\} \text{ and we assume}$$

$$(2.4) \quad \theta_n(\omega) \uparrow +\infty \quad \forall \omega.$$

For any $v \in V$, $\underline{x} \in E$, we define, as in [11], the following sequence of probability measures on (Ω, F_∞) , if $\xi_0 = i$

$$P^0 = P_{i,x}$$

P^1 is the (unique) probability measure on (Ω, F_∞) such that

$$P^1 = P^0 \text{ on } F_{\theta_1}$$

$$P^1(\eta_{\theta_1} B | F_{\theta_1}) = P_{\xi_1, \underline{x}_{\theta_1}}(B), \quad P^0 \text{ a.s.,}$$

$$\forall B \text{ Borel subset of } \Omega, \eta_t \text{ being the shift operator,}$$

and so on \cdots

$$P^n \text{ is similarly defined from } P^{n-1}$$

$$P^n = P^{n-1} \text{ on } F_{\theta_n}$$

$$P^n(\eta_{\theta_n} B | F_{\theta_n}) = P_{\xi_n, \underline{x}_{\theta_n}}(B), \quad P^{n-1} \text{ a.s.}$$

Defining

$$\xi(t) = \xi_n \quad \text{for } t \in [\theta_n, \theta_{n+1}[, \quad n \geq 0.$$

We consider the discounted reward

$$J_x(v) = \lim_{n \uparrow \infty} E_x^n \int_0^{\theta_{n+1}} e^{-\alpha t} f(\xi(t), \underline{x}_t) dt$$

where

$$f(\xi, \underline{x}) = f_i(x_i) \quad \text{iff } \xi = i.$$

Actually, with the assumptions $\theta_n \uparrow +\infty$, one can know that there exists a unique probability measure $P_{i,\underline{x}}^v$ on (Ω, F_∞) such that

$$(2.5) \quad P_{i,\underline{x}}^v = P^n \quad \text{on } F_{\theta_n}$$

and one can also define our total reward by

$$(2.6) \quad J_{\underline{x}}(v) = E_{\underline{x}}^v \int_0^\infty e^{-\alpha t} f(\xi_t, \underline{x}_t) dt.$$

We now add another control possibility, namely the “retirement option.”

Let T be the set of F_t stopping times, for $v \in V$, $\tau \in T$, and $(i, \underline{x}) \in U \times E$, we define the total reward as

$$(2.7) \quad J_{\underline{x}}^M(v, \tau) = E_{\underline{x}}^v \left\{ \int_0^\tau e^{-\alpha t} f(\xi_t, \underline{x}_t) dt + e^{-\alpha \tau} M \right\}$$

where M is a given constant.

We will use, as in Whittle [9], the additional assumption

$$(2.8) \quad \alpha k \leq f_j \leq \alpha K, \quad \forall j \in U,$$

where $k < K$ are given nonnegative constants.

The optimal reward function is

$$(2.9) \quad u(\underline{x}, M) = \text{Sup} (J_{\underline{x}}^M(v, \tau), (v, \tau) \in V \times T).$$

Using a *formal* dynamic programming argument, it is easy to check that $u(\underline{x}, M)$ should solve the following inequalities

$$(2.10) \quad \begin{aligned} w(\underline{x}, M) &\geq e^{-\alpha t} \Phi^i(t) w + \int_0^t e^{-\alpha s} \Phi^i(s) f_i(x_i) ds, \quad \forall t > 0 \quad \forall i \in U, \\ w(\underline{x}, M) &\geq M, \\ w(\cdot, M) &\text{ is a bounded measurable function.} \end{aligned}$$

In the following section, we will show that u is actually the *minimum* solution of these inequalities (for fixed M).

Let us recall the following result (cf. Bensoussan and Robin [3], Bensoussan [1]):

THEOREM 2.1. *Under the assumption (2.2) there exists a minimum solution \bar{u} of (2.10) in the space of bounded measurable functions. Moreover \bar{u} is upper semicontinuous.*

Remark 2.1. In Bensoussan and Robin [3], another kind of interpretation was given for $\bar{u}(\underline{x}, M)$. The present one will be more suitable for the problem we consider.

2.2. Characterization of the optimal reward (2.9). In order to characterize $u(\underline{x}, M)$ as defined in (2.9), we introduce another switching problem, with a switching cost ε . Namely, we consider the same problem as in § 2.1, but now, at each switching time a cost ε (i.e., a reward $-\varepsilon$) is involved. This is in fact a classical switching problem (which can be considered as an impulse control problem where the state is (ξ_t, \underline{x}_t) , cf. Bensoussan [1], Bensoussan and Lions [2] for the general theory).

In this context, let

$$V_0 = \{v = (\theta_n, \xi_n)_{n \geq 1}\}$$

be the set of admissible controls, θ_n, ξ_n being defined as previously.

For $(i, \underline{x}) \in U \times E$, define the reward

$$(2.11) \quad J_{i, \underline{x}}^{M, \varepsilon}(v, \tau) = E_{i, \underline{x}}^v \left\{ \int_0^\tau e^{-\alpha s} f(\xi(s), \underline{x}_s) dt - \varepsilon \sum_{j \neq i} e^{-\alpha \theta_j} \chi_{\theta_j < \tau} + e^{-\alpha \tau} M \right\}$$

where $E_{i, \underline{x}}^v$ is defined as in § 2.1, $(v, \tau) \in V_0 \times T$, and $\chi_B(\omega)$ is the characteristic function of the set B and $\xi_0 = i$ for the construction of $P_{i, \underline{x}}^n$.

We also define

$$(2.12) \quad u_i^\varepsilon(\underline{x}, M) = \sup (J_{i, \underline{x}}^{M, \varepsilon}(v, \tau), (v, \tau) \in V_0 \times T).$$

Let $u^\varepsilon = (u_1^\varepsilon, \dots, u_N^\varepsilon)$.

From impulse control theory (cf. Bensoussan and Lions [2], [11]) we know that, for fixed M , u^ε is the minimum element of the set of bounded measurable functions w satisfying

$$(2.13) \quad \begin{aligned} w_i(\underline{x}) &\geq e^{-\alpha t} \Phi^i(t) w_i + \int_0^t e^{-\alpha s} \Phi^i(s) f_i(x_s) ds, \quad \forall t > 0, \\ w_i(\underline{x}) &\geq -\varepsilon + \max_j w_j(\underline{x}), \\ w_i(\underline{x}) &\geq M. \end{aligned}$$

Moreover, $u_i^\varepsilon(\underline{x}) \in C(E)$, $\forall i = 1, \dots, N$.

We first establish the following result.

THEOREM 2.2. *Let $\underline{u}(\underline{x}, M)$ be the minimum solution of the inequalities (2.10), then*

$$(2.14) \quad \lim_{\varepsilon \downarrow 0} u_i^\varepsilon(\underline{x}, M) = \underline{u}(\underline{x}, M)$$

pointwise in \underline{x} .

Proof. It is clear that $u_i^\varepsilon(\underline{x}, M)$ increases when ε decreases, and that $u_i^\varepsilon(\underline{x}, M)$ is bounded (say by $(1/\alpha)\|f\| + M$). Let us define

$$\underline{w}_i = \lim_{\varepsilon \downarrow 0} u_i^\varepsilon.$$

From (2.13), we have

$$(2.15) \quad \underline{w}_i \geq \max_j w_j(\underline{x}, M), \quad \forall i.$$

Hence

$$\underline{w}_i(\underline{x}, M) = \underline{w}(\underline{x}, M) \quad \forall i.$$

But $\underline{u}(\underline{x}, M)$, the minimum solution of (2.10), satisfies obviously (2.13) and therefore

$$u_i^\varepsilon(\underline{x}, M) \leq \underline{u}(\underline{x}, M) \quad \forall \varepsilon, i.$$

So we deduce, when $\varepsilon \rightarrow 0$,

$$(2.16) \quad \underline{w}(\underline{x}, M) \leq \underline{u}(\underline{x}, M).$$

But we see that $\underline{w}(\underline{x}, M)$ will also satisfy (2.10), since this is identical to (2.13) when $\varepsilon = 0$ for a function which does not depend explicitly on i .

Therefore

$$\underline{w}(\underline{x}, M) \geq \underline{u}(\underline{x}, M).$$

Hence

$$\underline{w}(\underline{x}, M) = \underline{u}(\underline{x}, M)$$

and the theorem is proved. \square

Let us define

$$(2.17) \quad u(\underline{x}, M) = \sup (J_x^M(v, \tau), (v, \tau) \in V \times T).$$

Then we have the following Theorem.

THEOREM 2.3.

$$(2.18) \quad \underline{u}(\underline{x}, M) = u(\underline{x}, M).$$

Proof. Since $\varepsilon > 0$, we have

$$u_i^\varepsilon(\underline{x}, M) = \sup (J_{i,\underline{x}}^{M,\varepsilon}(v, \tau), (v, \tau) \in V_0 \times T)$$

and

$$J_{i,\underline{x}}^{M,\varepsilon}(v, \tau) \leq J_x^M(\tilde{v}, \tau)$$

where $\tilde{v} = ((i, 0), v)$, and $(\tilde{v}, \tau) \in V \times T$.

Therefore

$$u_i^\varepsilon(\underline{x}, M) \leq u(\underline{x}, M),$$

hence

$$(2.19) \quad \underline{u}(\underline{x}, M) \leq u(\underline{x}, M).$$

Now, for any solution w of the inequalities (2.10), one can show as in [11, Thm. VII, § 3.1] or [2b, § 6.4], that

$$w(\underline{x}) \geq E_x^m \left\{ e^{-\alpha \theta_{m+1} \wedge \tau} w(\underline{x}_{\theta_{m+1} \wedge \tau}) + \int_0^{\theta_{m+1} \wedge \tau} e^{-\alpha t} f(\underline{x}_t, v_t) dt \right\}$$

for any admissible control (v, τ) , $v = (\theta_i, \xi_i)_{i \geq 0}$ where E_x^m is the expectation corresponding to the measure P_x^m associate to (v, τ) as in § 1, with $\theta_m \wedge \tau$ instead of θ_m .

From this inequality, we deduce, when $m \rightarrow +\infty$, since $w(\underline{x}) \geq M$ and $\theta_m \wedge \tau \uparrow \tau$, that

$$w(\underline{x}) \geq J_x^M(v, \tau)$$

and therefore

$$w(\underline{x}) \geq u(\underline{x}, M).$$

Finally, this gives

$$\underline{u}(\underline{x}, M) \geq u(\underline{x}, M),$$

which, with (2.19), proves the result. \square

2.3. Reduction to write off policies. Following Whittle, a write off policy is defined as a policy such that there exists a family of “write-off” sets $S_i \subset E_i$ with the following properties.

- as soon as x_i (the state of the process i) belongs to S_i , the process i is abandoned;
- one retires as soon as all the processes have been abandoned, and only then;
- before retiring, one works only with those processes which have not been abandoned.

In this section we are going to show a lemma similar to the one obtained by Tsitsiklis [6] for discrete time, showing that we can restrict ourselves to write off policies with write off sets defined by optimal stopping problems for the individual processes.

The individual stopping problems. Let us consider the optimal stopping reward

$$(2.20) \quad \phi_i(x_i, M) = \sup_{\tau} I_{x_i}^M(\tau)$$

$$(2.21) \quad I_{x_i}^M(\tau) = E_{x_i}^i \left[\int_0^{\tau} e^{-\alpha t} f_i(x_{it}) dt + e^{-\alpha \tau} M \right].$$

It is known from standard theory (see Bensoussan [1]) that $\phi_i(x_i, M)$ is the minimum element of the set of functions $w(x_i)$ satisfying

$$(2.22) \quad w(x_i) \geq e^{-\alpha t} \Phi^i(t) w + \int_0^t e^{-\alpha s} \Phi^i(s) f_i(x_i) ds, \quad \forall t > 0$$

$$w(x_i) \geq M, \quad w \in C(E_i).$$

Let us show the following results which extend the discrete time case (cf. Whittle [9]) and the diffusion case (Karatzas [5]).

LEMMA 2.1. *Under the assumptions (2.2)-(2.8), $\phi(x, M) = \phi_i(x_i, M)$ has the following properties:*

- (i) $\phi(x, M) = M \quad \forall M \geq K$;
- (ii) $\phi(x, M) = \int_0^{\infty} e^{-\alpha t} \Phi^i(t) f_i(x) dt, \quad \forall M \leq k$;
- (iii) $\forall x \in E_i, \phi(x, \cdot)$ is an increasing convex function;
- (iv) $\phi(x, \cdot)$ is Lipschitz continuous and in every M where the derivative exists

$$0 \leq \frac{\partial \phi}{\partial M} \leq 1;$$

- (v) in every point where the derivative exists

$$\frac{\partial \phi}{\partial M}(x, M) = E_x e^{-\alpha \hat{\tau}}$$

where $\hat{\tau}$ is optimal for (2.20), namely

$$\hat{\tau} = \inf(t \geq 0, \phi(x, M) = M).$$

Proof. (i) This shows that $\tau = 0$ is optimal in (2.20) whenever $M \geq K$. Since $f_i \leq \alpha K$

$$J_x^M(\tau) \leq E_x \{ (1 - e^{-\alpha \tau}) K + e^{-\alpha \tau} M \}$$

$$= K + E_x e^{-\alpha \tau} (M - K),$$

clearly if $M \geq K, \tau = 0$ gives the maximum value.

- (ii) $f_i \geq \alpha k$ implies

$$w_0(x) = \int_0^{\infty} e^{-\alpha t} \Phi^i(t) f_i dt \geq k$$

and since $w_0(x) = e^{-\alpha t} \Phi^i(t) w_0 + \int_0^t e^{-\alpha s} \Phi^i(s) f_i ds$, we see that w_0 satisfies (2.22) for $M = k$.

Moreover $t = +\infty$ in (2.22) gives

$$w(x) \geq w_0(x) \quad \forall w \text{ solution of (2.22).}$$

- (iii) If $0 \leq \lambda \leq 1$, then we check that

$$w_{\lambda} = \lambda \phi(x, m_1) + (1 - \lambda) \phi(x, m_2)$$

satisfies (2.22) for $m = \lambda m_1 + (1 - \lambda)m_2$ and therefore

$$w_\lambda \geq \phi(x, \lambda m_1 + (1 - \lambda)m_2).$$

The increasing property is obvious from (2.20).

(iv) From (2.21) one has, for an arbitrary τ

$$I_x^{M+\delta}(\tau) - I_x^M(\tau) = E_x e^{-\alpha\tau}\delta$$

therefore, for $\delta > 0$

$$I_x^{M+\delta}(\tau) \leq I_x^M(\tau) + \delta \leq \phi(x, M) + \delta$$

implying

$$\phi(x, M + \delta) \leq \phi(x, M) + \delta$$

and since $\phi(x, M + \delta) \geq \phi(x, M)$, we see that

$$0 \leq \frac{\partial^+ \phi}{\partial M}(x, M) \leq 1.$$

(v) Let $\hat{\tau} = \inf(t \geq 0, \phi(x, M) = M)$, we know (cf. [1]) that

$$\phi(x, M) = I_x^M(\hat{\tau}).$$

Therefore, if $\delta > 0$,

$$\begin{aligned} I_x^{M+\delta}(\hat{\tau}) &= I_x^M(\hat{\tau}) + \delta E_x e^{-\alpha\hat{\tau}} \\ &= \phi(x, M) + \delta E_x e^{-\alpha\hat{\tau}} \end{aligned}$$

hence

$$\begin{aligned} \phi(x, M + \delta) - \phi(x, M) &\geq \delta E_x e^{-\alpha\hat{\tau}} \\ \frac{\partial^+ \phi}{\partial M}(x, M) &\geq E_x e^{-\alpha\hat{\tau}}. \end{aligned}$$

Taking $\delta < 0$, we get

$$\frac{\partial^- \phi}{\partial M}(x, M) \leq E_x e^{-\alpha\hat{\tau}}.$$

Therefore, in M such that the derivative exists, we get the result. \square

COROLLARY. $\partial^+ \phi / \partial M$ is a right continuous increasing function such that

$$\begin{aligned} 0 &\leq \frac{\partial^+ \phi}{\partial M} \leq 1, \\ \frac{\partial^+ \phi}{\partial M}(x, M) &= 1 \quad \forall M \geq K \\ \frac{\partial^+ \phi}{\partial M}(x, M) &= 0 \quad \forall M < k. \end{aligned}$$

Let us now define, for fixed i

$$y_i = (x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_N)$$

$U_i(y_i, M)$ the optimal reward function when only the processes different from i are available.

From the previous section, $U_i(y_i, M)$ is the minimum solution of

$$(2.23) \quad W_i(y_i, M) \geq e^{-\alpha t} \Phi^j(t) W_i + \int_0^t e^{-\alpha s} \Phi^j(s) f_j(x_j) ds, \quad \forall t > 0, \quad \forall j \neq i$$

$$W_i(y_i, M) \geq M, \quad W_i(\cdot, M) \text{ bounded and measurable.}$$

We can now state the Tsitsiklis' lemma in continuous time.

LEMMA 2.2. (*Tsitsiklis' lemma in continuous time*). *One has*

$$(2.24) \quad u(\underline{x}, M) \leq \phi_i(x_i, M) - M + U_i(y_i, M), \quad \forall i.$$

Proof. Let $W_i(\underline{x}, M)$ be the right-hand side of (2.24). We are going to show that W_i satisfies (2.10) and since u is the minimum solution, this will show the lemma. Notice that since $U_i \geq M$ and $\phi_i \geq M$, we have

$$W_i(\underline{x}, M) \geq M.$$

Moreover, since U_i does not depend on x_i ,

$$\begin{aligned} e^{-\alpha t} \Phi^i(t) W_i + \int_0^t e^{-\alpha s} \Phi^i(s) f_i ds &= \left\{ e^{-\alpha t} \Phi^i(t) \phi_i + \int_0^t e^{-\alpha s} \Phi^i(s) f_i ds \right\} + e^{-\alpha t} [U_i - M] \\ &= I + II. \end{aligned}$$

We have,

$$I \leq \phi_i \quad \text{by (2.22)}$$

and, since $U_i - M \geq 0$,

$$II \leq U_i - M.$$

Then, for $j \neq i$, since ϕ_i does not depend on x_j , $j \neq i$,

$$\begin{aligned} e^{-\alpha t} \Phi^j(t) W_i + \int_0^t e^{-\alpha s} \Phi^j(s) f_j ds \\ = \left\{ e^{-\alpha t} \Phi^j(t) U_i + \int_0^t e^{-\alpha s} \Phi^j(s) f_j ds \right\} + e^{-\alpha t} [\phi_i - M] = III + IV. \end{aligned}$$

Hence, using (2.23)

$$III \leq U_i,$$

and

$$IV \leq \phi_i - M \text{ since } \phi_i - M \geq 0.$$

Therefore the lemma is proved. \square

COROLLARY. Define $S_i^M = \{x_i \in E_i, \phi_i(x_i, M) = M\}$, then one can restrict the policies to be write off with respect to $(S_i^M, i = 1, \dots, N)$.

Proof. Notice that $U_i(y_i, M) \leq u(\underline{x}, M)$. If $x_i \in S_i^M$, then (2.24) gives, with the above inequality,

$$u(\underline{x}, M) = U_i(y_i, M)$$

which means that the optimal reward is the same as the one with $N-1$ processes where the i process has been dropped. If there exists i such that, for $\underline{x} = (x_1, \dots, x_N)$, $x_i \notin S_i^M$, then $u(\underline{x}, M) \geq \phi_i(x_i, M) > M$ implies that it is not optimal to retire. Finally

if \underline{x} is such that $x_i \in S_i^M, \forall i$, then $u(\underline{x}, M) = U_i \forall i$ and we can use the same argument for $N-1$ processes to show that

$$u(\underline{x}, M) = \phi_j(x_j, M) = M \quad \forall j. \quad \square$$

Let us denote by V_{off}^M the set of admissible write off policies corresponding to $(S_i^M, i = 1, \dots, N)$. We will use the following lemma due to Whittle (cf. [9]).

LEMMA 2.3. *If (v, τ) is a write off policy, then*

$$E_{\underline{x}}^v e^{-\alpha\tau} = \prod_{i=1}^N E_{x_i} e^{-\alpha\tau_i}$$

where τ_i is the retirement times when only the process i is available.

Proof. For the proof see Whittle [9]. \square

In our context, this means that, if

$$(2.25) \quad \tau_i^M = \inf (t \geq 0, \phi_i(x_{it}, M) = M),$$

then, for all write off policies,

$$(2.26) \quad E_{\underline{x}}^v e^{-\alpha\tau} = \prod_i E_{x_i} e^{-\alpha\tau_i^M}.$$

We can then deduce the product formula for u .

THEOREM 2.4.

$$(2.27) \quad U(\underline{x}, M) = K - \int_M^K \prod_{i=1}^N \frac{\partial \phi_i}{\partial m}(x_i, m) dm.$$

Proof. Let (v, τ) be an admissible write off policy with respect to $(S_i^M, i = 1, \dots, N)$. We have,

$$J_{\underline{x}}^{M+\delta}(v, \tau) - J_{\underline{x}}^M(v, \tau) = \delta E_{\underline{x}}^v e^{-\alpha\tau}.$$

From the previous lemma

$$J_{\underline{x}}^{M+\delta}(v, \tau) - J_{\underline{x}}^M(v, \tau) = \delta \cdot \prod_{i=1}^N E_{x_i} e^{-\alpha\tau_i^M},$$

therefore

$$u(\underline{x}, M + \delta) \geq J_{\underline{x}}^M(v, \tau) + \delta \prod_{i=1}^N E_{x_i} e^{-\alpha\tau_i^M}.$$

Note that the last term is independent from (v, τ) as far as (v, τ) is a write off policy with respect to (S_i^M) . Therefore, maximizing with respect to (v, τ)

$$u(\underline{x}, M + \delta) \geq u(\underline{x}, M) + \delta \prod_i E_{x_i} e^{-\alpha\tau_i^M}$$

which implies, for $\delta > 0$,

$$\frac{\partial^+ u}{\partial M}(\underline{x}, M) \geq \prod_i E_{x_i} e^{-\alpha\tau_i^M}$$

and for $\delta < 0$

$$\frac{\partial^- u}{\partial M}(\underline{x}, M) \leq \prod_i E_{x_i} e^{-\alpha\tau_i^M}.$$

Therefore, at every point where the derivatives exist, we have, thanks to the Lemma 2.1,

$$\frac{\partial u}{\partial M}(x, M) = \prod_{i=1}^N \frac{\partial \phi_i}{\partial M}(x_i, M).$$

Integrating from M to K , using the fact that $u(x, K) = K$, we get (2.27). \square

Remark. From Bensoussan and Robin [3], we can show that the optimal reward of the discrete time problem converges to the $u(x, M)$ when the time step h goes to zero. However, we have not been able to show the product formula in continuous time by letting h go to zero on the product formula of the discrete time case.

2.4. The forward induction lemma. Let us consider the discrete time version of the stopping problems (2.22). Namely, for $h > 0$, we define (dropping the index i)

$$\begin{aligned} r_h(x) &= E_x \int_0^h e^{-\alpha s} f(x_s) ds \\ Q_h z &= \phi(h)z \\ \beta &= e^{-\alpha h}. \end{aligned}$$

Then the optimal reward for the discrete stopping problem $\phi_h(x, m)$ is the unique solution of

$$\phi_h(x, m) = \max (r_h + \beta Q_h \phi_h, m).$$

Defining $V_h = \{\tau, \text{stopping times with values in } N_h = \{nh, n \geq 0\}\}$, we can write

$$\phi_h(x, m) = \sup_{\tau \in V_h} E_x \left(\int_0^\tau e^{-\alpha s} f(x_s) ds + e^{-\alpha \tau} m \right).$$

The *index* is defined, as previously, as

$$M_h(x) = \inf (m > k, \phi_h(x, m) = m).$$

On the other hand, Whittle [9] shows that $M_h(x)$ has the following representation:

$$(2.28) \quad M_h(x) = \sup_{\tau \in V_h^*} \frac{E_x \int_0^\tau e^{-\alpha s} f(x_s) ds}{1 - E_x e^{-\alpha \tau}}$$

with $V_h^* = \{\tau \text{ stopping times with values in } N_h^* = N_h - \{0\}\}$.

The extension of the formula (2.28) to diffusion processes was done by Karatzas [5] using explicit calculation for one-dimensional processes. We are going to show the same formula in our context; the idea being to approximate the stopping problem (2.22) by a discrete time problem (like in Bensoussan-Robin [3]).

LEMMA 2.4. *Let $\phi(x, m)$ be defined as in (2.20) (where we drop the index i), and define*

$$M(x) = \inf (m > k, \phi(x, m) = m), \text{ and } V^* = \bigcup_h V_h^*,$$

then

$$M(x) = \sup_{\tau \in V^*} \frac{E_x \int_0^\tau e^{-\alpha s} f(x_s) ds}{1 - E_x e^{-\alpha \tau}}$$

where $V^* = \{\bigcup_h V_h^*\}$.

Proof. Starting with $M_h(x)$ we have

$$k \leq M_h(x) \leq K \quad \forall x, \quad \forall h.$$

Clearly, V_h^* is increasing as h decreases to zero and therefore $M_h(x)$ is increasing when $h \downarrow 0$. For fixed x , let

$$v(x) = \lim_{h \downarrow 0} M_h(x)$$

then

$$v(x) = \sup_{\tau \in V^*} Z(\tau), \quad \text{where } Z(\tau) = \frac{E_x \int_0^\tau e^{-\alpha s} f(x_s) ds}{1 - E_x e^{-\alpha \tau}}.$$

Indeed, for all ε , there exists h , such that

$$v \cong M_h > v - \varepsilon \Rightarrow \exists \delta(\varepsilon) \text{ s.t. } M_h - \delta(\varepsilon) > v - \varepsilon$$

and from the definition of M_h , we can find $\tau_h(\delta(\varepsilon))$ such that, $\tau_h \in V_h^*$,

$$M_h \cong Z(\tau_h) > M_h - \delta(\varepsilon).$$

Therefore for all ε , there exists $\tau \in V^*$ such that $v \cong Z(\tau) > v - \varepsilon$ proving that $v = \sup (Z(\tau), \tau \in V^*)$.

Let us prove that $v = M(x)$. Assume that $m \cong v$, then

$$m \cong v \cong \frac{E_x \int_0^\tau e^{-\alpha s} f(x_s) ds}{1 - E_x e^{-\alpha \tau}} \quad \forall \tau \in V^*$$

\Rightarrow

$$m \cong E_x \left(\int_0^\tau e^{-\alpha s} f(x_s) + e^{-\alpha \tau} m \right) \quad \forall \tau \in V^*$$

(and for $\tau = 0$, we have the equality). Therefore $m \cong \phi(x, m)$.

But $\phi(x, m) \cong m$, for all m implies $\phi(x, m) = m$ for all $m \cong v$. Now assume that $m < v$.

Let us assume that for such m

$$\phi(x, m) = \sup_{\tau} E_x \int_0^\tau e^{-\alpha s} f(x_s) + e^{-\alpha \tau} m = m.$$

This would imply

$$m \cong \frac{E_x \int_0^\tau e^{-\alpha s} f(x_s) ds}{1 - E_x e^{-\alpha \tau}} \quad \forall \tau \in V^*$$

which contradicts the assumption $m < v$.

Therefore

$$m < v \Rightarrow \phi(x, m) > m, \quad \text{hence } v = M(x). \quad \square$$

Remark. As it was stressed in Katehakis-Veinott [12] we can also characterize $M_h(x)$ using the “restart in x -problem” for which the optimal reward function $v_h^x(\cdot)$ is given by

$$v_h^x(y) = \max (r_h(y) + \beta Q_h v_h^x, r_h(x) + \beta Q_h v_h)$$

and then (see [12]) we have

$$M_h(x) = v_h^x(x).$$

In continuous time, in order to define a similar problem, we can use the discrete time solution; namely, as in Bensoussan–Robin [3], we could show that for $h = 2^{-N}$

$$v_N^x(y) = v_h^x(y)$$

is increasing when $N \rightarrow +\infty$, (and bounded).

Then

$$v^x(y) = \lim_{N \uparrow \infty} v_N^x(y)$$

is the minimum solution of the inequalities

$$v^x(y) \geq e^{-\alpha t} \phi(t) v^x + \int_0^t e^{-\alpha s} \Phi(s) f ds$$

$$v^x(y) \geq v^x(x), \quad \forall y$$

$v^x(x)$ bounded measurable functions.

This is the continuous time version of the restart in x -problem.

3. The problem with switching cost. We now turn back to the case where there is a switching cost incurred at each time we change the active process. This was already considered in § 2.2 when we constructed the functions u_i^ε . Recall that this is a more or less standard impulse control problem where the underlying state is in fact (z, \underline{x}) where $z \in \{1, \dots, N\}$ is the number of the active process. It would be interesting to know if a product formula like (1.4) holds. We do not know the answer, neither for the question of the optimality of some index rule. However, we can show that the concept of write off policy is still valid in this case and this gives some more information on the optimal policies than the mere interpretation of the dynamic programming condition. The reduction to write off policies will be a consequence of the following simple result, similar to the Tsitsiklis' lemma. Let us make precise some notations: we drop the ε in the optimal reward which is now

$$(3.1) \quad u(z, \underline{x}, M) = \sup (J_{z, \underline{x}}^M(v, \tau), (v, \tau) \in V_0 \times T)$$

$J_{z, \underline{x}}^M(v, \tau)$, V_0 , T being defined as in (2.11), with $z \in \{1, \dots, N\}$. We know that u is the minimum element of the set of bounded and measurable functions $w(z, \underline{x})$ satisfying

$$w(z, \underline{x}) \geq e^{-\alpha t} \Phi^z(t) w + \int_0^t e^{-\alpha s} \Phi^z(s) f_z(x_z) ds$$

$$(3.2) \quad w(z, \underline{x}) \geq -\varepsilon + \max_j w(j, \underline{x}),$$

$$w(z, \underline{x}) \geq M, \quad \forall z \in \{1, \dots, N\}.$$

We denote by

$$y_i = (x_j, j \neq i),$$

$U(z, y_i, M)$ the optimal reward when only the processes different from i are available, and when the initial active process is the process number z .

LEMMA 3.1. *We have for arbitrary $i \in \{1, \dots, N\}$,*

$$(3.3) \quad u(j, \underline{x}, M) \leq [\phi_i(x_i, M) - (M + \varepsilon)]^+ + U(j, y_i, M) \quad \forall j \neq i$$

$$(3.4) \quad u(i, \underline{x}, M) \leq \phi_i(x_i, M) - M + \max_{j \neq i} \left[M, -\varepsilon + \max_{j \neq i} U(j, y_i, M) \right].$$

Proof. Let us define, for fixed i :

$$w(z, \underline{x}) = \begin{cases} [\phi_i(x_i, M) - (M + \varepsilon)]^+ + U(z, y_i, M) & \text{for } z \neq i \\ \phi_i(x_i, M) - M + \max \left[M, -\varepsilon + \max_{j \neq i} U(j, y_i, M) \right] & \text{for } z = i. \end{cases}$$

We are going to show that $w(z, \underline{x})$ satisfies (3.2) and since $u(z, \underline{x}, M)$ is the minimum solution, this will prove the lemma. We have,

$$\begin{aligned} e^{-\alpha t} \Phi^z(t) w(z, \underline{x}) + \int_0^t e^{-\alpha s} \Phi^z(s) f_z(x_z) ds \\ = \left\{ e^{-\alpha t} \Phi^z(t) U(z, y_i, M) + \int_0^t e^{-\alpha s} \Phi^z(s) f_z(x_z) ds \right\} \\ + e^{-\alpha t} [\phi_i(x_i, M) - (M + \varepsilon)]^+ \quad \text{if } z \neq i \\ = \left\{ e^{-\alpha t} \Phi^i(t) \phi_i + \int_0^t e^{-\alpha s} \Phi^i(s) f_i(x_i) ds \right\} \\ + e^{-\alpha t} \left[\max \left[M, -\varepsilon + \max_{j \neq i} U(j, y_i, M) \right] - M \right] \quad \text{if } z = i. \end{aligned}$$

In the first case, thanks to (3.2), the right-hand side is less than

$$U(z, y_i, M) + e^{-\alpha t} [\phi_i(x_i, M) - (M + \varepsilon)]^+$$

i.e., less than

$$U(z, y_i, M) + [\phi_i(x_i, M) - (M + \varepsilon)]^+ = w(z, \underline{x}).$$

In the second one, thanks to (2.22) for ϕ_i , the right-hand side is less than

$$\phi_i(x_i, M) - M + \max \left[M, -\varepsilon + \max_{j \neq i} U(j, y_i, M) \right] = w(z, \underline{x}) \quad \text{if } z = i.$$

Therefore the first inequality of (3.2) is satisfied. It is obvious that $w(z, \underline{x}) \geq M$. Now, for the second inequality of (3.2), we must check that

$$(3.5) \quad w(i, \underline{x}) \geq -\varepsilon + \max_{j \neq i} ([\phi_i(x_i, M) - (M + \varepsilon)]^+ + U(j, y_i, M))$$

and, for $z \neq i$

$$(3.6) \quad w(z, \underline{x}) \geq -\varepsilon + \max \left\{ \begin{aligned} & [\phi_i(x_i, M) - (M + \varepsilon)]^+ + U(j, y_i, M) \quad \forall j \neq i \\ & \phi_i(x_i, M) - M + \max \left[M, -\varepsilon + \max_{j \neq i} U(j, y_i, M) \right]. \end{aligned} \right.$$

But, since $\phi_i(x_i, M) - M \geq [\phi_i(x_i, M) - (M + \varepsilon)]^+$, (3.5) is obvious from the definition of $w(i, \underline{x})$. For (3.6), since

$$U(z, y_i, M) \geq -\varepsilon + \max_{j \neq i} U(j, y_i, M)$$

we have

$$w(z, \underline{x}) \geq -\varepsilon + [\phi_i(x_i, M) - (M + \varepsilon)]^+ + \max_{j \neq i} U(j, y_i, M)$$

and since $[\phi_i(x_i, M) - (M + \varepsilon)]^+ \geq \phi_i(x_i, M) - M - \varepsilon$, we also have

$$w(z, \underline{x}) \geq -\varepsilon + \phi_i(x_i, M) - M + \max \left[M, -\varepsilon + \max_{j \neq i} U(j, y_i, M) \right].$$

Therefore, the lemma is obtained. \square

Let us define the following write off sets:

$$(3.7) \quad S_i^M = \{(z, \underline{x}) \in \{1, \dots, N\} \times E_i \text{ such that either } z = i \text{ and } \phi_i(x_i, M) = M, \\ \text{or } z \neq i \text{ and } \phi_i(x_i, M) \leq M + \varepsilon\}.$$

THEOREM 3.1. *We can restrict the admissible policies to be write off with respect to $(S_i^M, i = 1, \dots, N)$, in other words*

- (i) *if $\exists i$ s.t. $(z, \underline{x}) \notin S_i^M$, we continue (i.e., we do not use the retirement option)*
- (ii) *if $\forall i, (z, \underline{x}) \in S_i^M$, we retire*
- (iii) *if $(z, \underline{x}) \in S_i^M$, the process i is abandoned.*

Proof. (i) Assume that $\exists i$ s.t. $(z, \underline{x}) \notin S_i^M$

then -either $z = i$ and $\phi_i(x_i, M) > M$ hence $u(i, \underline{x}, M) \geq \phi_i(x_i, M) > M$,

therefore we do not retire;

-or $z \neq i$ and $\phi_i(x_i, M) > M + \varepsilon$

hence $u(z, \underline{x}, M) \geq -\varepsilon + \max_j u(j, \underline{x}, M) \geq -\varepsilon + \phi_i > M$.

Therefore we do not retire.

(ii) Assume that

$$(3.8) \quad \forall i, (z, \underline{x}) \in S_i^M$$

and to fix the idea, take $z = N$, then (3.3) implies, since $(z, \underline{x}) \in S_1^M$, and $U(z, y_1, M) \leq u(z, \underline{x}, M)$,

$$u(z, \underline{x}, M) = U(z, y_1, M).$$

Denote U by $U^{N-1}(z, y_1^{N-1}, M)$ to make explicit that U is the optimal reward of a problem where only the $N - 1$ first components are available, i.e., $y_1^{N-1} = (x_2, \dots, x_N)$.

Then applying again (3.3) to the $N - 1$ dimensional bandit problem we get, with $i = 2$

$$u(z, \underline{x}, M) = U^{N-1}(z, y_1^{N-1}, M) = U^{N-2}(z, y_2^{N-2}, M)$$

with $y_2^{N-1} = (x_3, \dots, x_N)$.

This process goes on until

$$u(z, \underline{x}, M) = U^1(z, x_z, M) = \phi_z(x_z, M)$$

which by the assumption and (3.4) is equal to M . Therefore we must retire if (3.8) holds.

(iii) Assume $(z, \underline{x}) \in S_i^M$

-either $z \neq i$ then (3.3) and $u(z, \underline{x}, M) \geq U(z, y_i, M)$ implies $u(z, \underline{x}, M) = U(z, y_i, M)$ meaning that we never use again the process i

-or $z = i$ and $\phi_i(x_i, M) = M$, then (3.4) implies that either we retire, or we have

$$u(z, \underline{x}, M) = -\varepsilon + \max_{j \neq z} U(j, y_z, M)$$

meaning that we switch to another process and never use the process $z = i$. This completes the proof of Theorem 3.1. \square

Acknowledgment. The authors would like to thank Professors A. Bensoussan and P. L. Chow for the useful discussions on this work.

REFERENCES

- [1] A. BENSOUSSAN, *Stochastic Control by Functional Analysis Methods*, North-Holland, Amsterdam, 1982.
- [2a] A. BENSOUSSAN AND J. L. LIONS, *Inéquations variationnelles et problèmes d'arrêt optimal*, Dunod, Paris, 1978.
- [2b] ———, *Applications des inéquations quasi variationnelles au contrôle stochastique*, Dunod, Paris, 1982.
- [3] A. BENSOUSSAN AND M. ROBIN, *On the convergence of the discrete time dynamic programming equation for general semi-group*, SIAM J. Control Optim., 20 (1982), pp. 722–746.
- [4] J. C. GITTINS, *Bandit processes and dynamic allocation indices*, J. Roy. Stat. Soc., 41 (1979), pp. 148–177.
- [5] I. KARATZAS, *Gittins indices in the dynamic allocation problem for diffusion processes*, Ann. Probab., 12 (1984), pp. 173–192.
- [6] J. N. TSITSIKLIS, *A lemma on the multiarmed bandit problem*, IEEE Trans. Automat. Control, AC-31 (1986), pp. 576–577.
- [7] P. VARAIYA, J. WALRAND, AND C. BUYUKKOC, *Extensions of the multiarmed bandit problem: the discounted case*, IEEE Trans. Automat. Control, AC-30 (1985), pp. 426–439.
- [8] P. WHITTLE, *Multiarmed bandit and the Gittins index*, J. Roy. Stat. Soc., 42 (1980), pp. 143–149.
- [9] ———, *Optimization Over Time*, vol. 1, John Wiley, New York, 1982.
- [10] E. B. DYNKIN, *Markov Processes*, vols. 1 and 2, Springer-Verlag, Berlin, New York, 1968.
- [11] M. ROBIN, *Contrôle impulsionnel des processus de Markov*, Thèse, Paris IX, 1978.
- [12] M. N. KAHETAKIS AND A. F. VEINOTT, *The multiarmed bandit problem: decomposition and computation*, Math. Oper. Res., 12 (1987), pp. 262–268.
- [13] A. MANDELBAUM, *Discrete multiarmed bandit and multi-parameter processes*, Probab. Theory Related Fields, 71 (1986), pp. 129–147.
- [14] ———, *Continuous multiarmed bandits and multi-parameter processes*, Ann. Probab., 15 (1987), pp. 1527–1556.